

Coordination-free Byzantine Replication with Minimal Communication Costs

Jelle Hellings Mohammad Sadoghi

Exploratory Systems Lab,
Department of Computer Science,
University of California, Davis, CA, USA



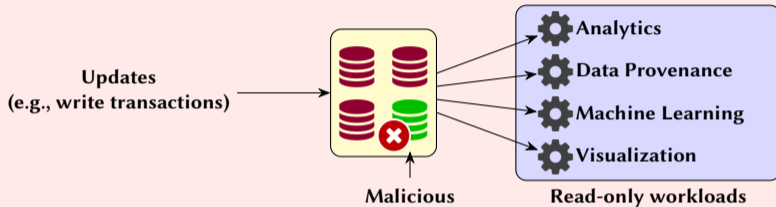
Towards high performance fault-tolerant data processing

Fault tolerance via full replication

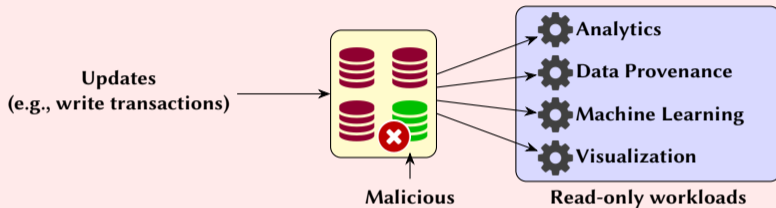
- ▶ *Communication intensive*: coordinate all steps of the system via *consensus*.
- ▶ *No scalability*: adding replicas reduces performance.
- ▶ *Replicas can be malicious*: read-only queries cannot be answered by single replicas.
- ▶ No sharding, no geo-aware data localization, no specialization.

Contradicts standard techniques used for high-performance data processing!

Towards our vision: specializing for read-only workloads



Towards our vision: specializing for read-only workloads



Enabling data-hungry read-only workloads

Need to stream all data updates with low cost for all replicas involved.

The need for Byzantine learning

Definition

Let \mathfrak{R} be a cluster deciding on a sequence of transactions.

The *Byzantine learning problem* is the problem of sending the decided transactions from \mathfrak{R} to a learner \mathcal{L} such that:

- ▶ the learner \mathcal{L} will eventually *receive all* decided transactions;
- ▶ the learner \mathcal{L} will *only receive* decided transactions.

Practical requirements

- ▶ Minimizing overall communication.
- ▶ Load balancing among all replicas in \mathfrak{R} .

Background: Information dispersal algorithms

Definition

Let v be a value with storage size $s = \|v\|$.

An *information dispersal algorithm* can encode v in \mathbf{n} pieces v' such that v can be *decoded* from every set of $\mathbf{n} - \mathbf{f}$ such pieces.

Theorem (Rabin 1989)

The IDA algorithm is an *optimal* information dispersal algorithm:

- ▶ Each piece v' has size $\left\lceil \frac{\|v\|}{\mathbf{n} - \mathbf{f}} \right\rceil$.
- ▶ The $\mathbf{n} - \mathbf{f}$ pieces necessary for decoding have a total size of $(\mathbf{n} - \mathbf{f}) \left\lceil \frac{\|v\|}{(\mathbf{n} - \mathbf{f})} \right\rceil \approx \|v\|$.

The delayed-replication algorithm

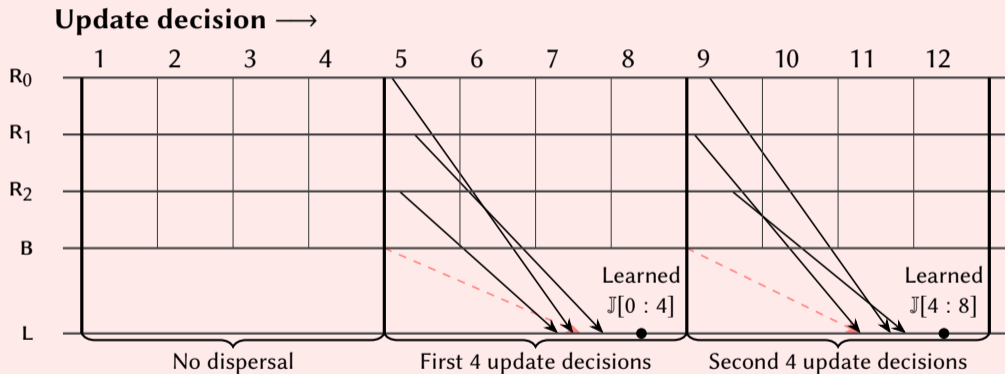
Idea: \mathfrak{R} sends a journal to learner \mathcal{L}

1. Partition the journal in sequences S of \mathbf{n} transactions.
2. Replica $\mathfrak{R}_i \in \mathfrak{R}$ encodes S into the i -th IDA piece S_i .
3. Replica $\mathfrak{R}_i \in \mathfrak{R}$ sends S_i with a checksum $C_i(S)$ of S to \mathcal{L} .
4. \mathcal{L} receives at least $\mathbf{n} - \mathbf{f}$ distinct pieces and decodes S .

Observation ($\mathbf{n} > 2\mathbf{f}$)

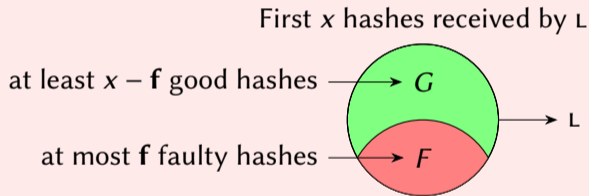
- ▶ Replica \mathfrak{R}_i sends at most $B = \left\lceil \frac{\|S\|}{\mathbf{n} - \mathbf{f}} \right\rceil + c \leq \frac{2\|S\|}{\mathbf{n}} + 1 + c = \mathcal{O}\left(\frac{\|S\|}{\mathbf{n}} + c\right)$ bytes.
- ▶ Learner \mathcal{L} receives at most $\mathbf{n} \cdot B = \mathcal{O}(\|S\| + \mathbf{c}\mathbf{n})$ bytes.

Communication by the delayed-replication algorithm



Decoding S using simple checksums ($n > 2f$)

- ▶ Use checksums $\text{hash}(S)$.
- ▶ The $n - f$ non-faulty replicas will provide correct *pieces*.
- ▶ At least $n - f > f$ messages with correct *checksums*.



Wait until $f + 1 \leq g$ identical hashes: $\text{hash}(S)$.

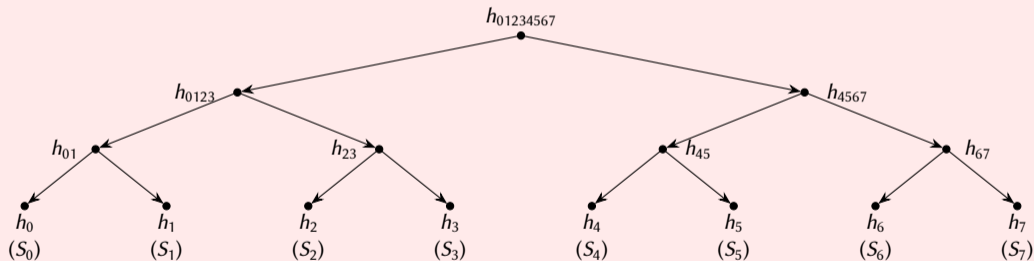
- ▶ Compute intensive for learners: one can choose $n - f$ out of n messages in $\binom{n}{n-f}$ ways *only one such choice is guaranteed to be correct!*

Decoding S using tree checksums

Use Merkle-trees to construct checksums

Consider 8 replicas and a sequence S .

We construct the checksum $C_5(S)$ of S (used by R_5).



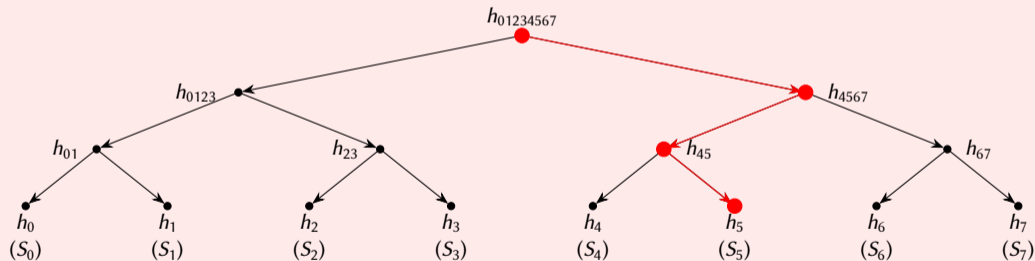
Construct a Merkle tree for pieces S_0, \dots, S_7 .

Decoding S using tree checksums

Use Merkle-trees to construct checksums

Consider 8 replicas and a sequence S .

We construct the checksum $C_5(S)$ of S (used by R_5).



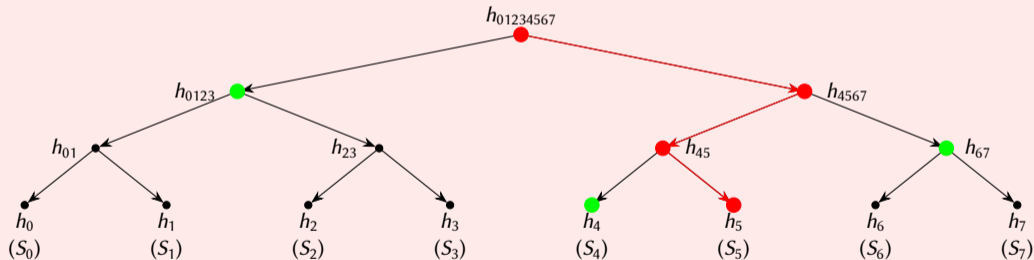
Determine the path from root to S_5 .

Decoding S using tree checksums

Use Merkle-trees to construct checksums

Consider 8 replicas and a sequence S .

We construct the checksum $C_5(S)$ of S (used by R_5).



Select *root* and *neighbors*: $C_5(S) = [h_4, h_{67}, h_{0123}, h_{01234567}]$.

Delayed-replication: Main result ($n > 2f$)

Theorem

Consider the learner L , replica R , and decided transactions \mathcal{T} . The delayed-replication algorithm with tree checksums guarantees

1. L will learn \mathcal{T} ;
2. L will receive at most $|\mathcal{T}|$ messages with a total size of $O(\|\mathcal{T}\| + |\mathcal{T}| \log n)$;
3. L will only need at most $|\mathcal{T}|/n$ decode steps;
4. R will send at most $|\mathcal{T}|/n$ messages to L of size $O\left(\frac{\|\mathcal{T}\| + |\mathcal{T}| \log n}{n}\right)$.

Application: scalable storage for resilient systems

- ▶ Clusters typically need only a *view* \mathbb{V} on the data to decide whether updates are valid.
- ▶ Clusters only need the full journal \mathbb{J} for *recovery*.
- ▶ We can use *delayed-replication* to reduce the data each replica has to store.

Theorem

The storage cost per replica can be reduced from

$$O(\|\mathbb{J}\| + \|\mathbb{V}\|) \quad \text{to} \quad O\left(\frac{\|\mathbb{J}\|}{\mathbf{n} - \mathbf{f}} + \frac{|\mathbb{J}|}{\mathbf{n}} \log(\mathbf{n}) + \|\mathbb{V}\|\right).$$

Conclusion

Efficient Byzantine learning is possible.

Blockchain applications

- ▶ Low-cost checkpoint protocols.
- ▶ Scalable storage for resilient systems.

Fault-tolerant high-performance data processing: ongoing work

- ▶ More at <https://jhellings.nl/> and <https://expolab.org/>.
- ▶ Paper: <https://doi.org/10.4230/LIPIcs.ICDT.2020.17>.